**Project 2 (15%)**
**STQD6014**
**SEMESTER 1 2023/2024**

As a skilled Data Scientist, your role involves uncovering and communicating the narratives hidden within datasets to stakeholders through visualizations. To initiate this process, please visit kaggle.com and select any **RAW DATA SET** of your preference, downloadable in **.CSV format** for the initial data exploration phase. Your tasks encompass:

1. Data Cleaning:
- Thoroughly clean the **RAW DATA** to ensure accuracy and reliability in subsequent analyses.

2. Visualizations:
- Generate a **minimum of five distinct visualization plots**, employing techniques such as pie charts, bar plots, box plots, etc.
- Utilize **Matplotlib**, **Seaborn** or other relevant packages of your choice for creating these visualizations.

3. Insights and Explanations:
- Accompany each plot with a **compelling insight and explanation** that is both relevant and significant. Provide clarity on the story the data tells. You are encouraged to substantiate your observations with other published information online.

4. **Google Colab Notebook** or **Jupyter Notebook**:
- Use the **Google Colab Notebook** or **Jupyter Notebook** as the platform for your assignment.
- Structure the Notebook like a book, encompassing sections such as:
    - **Introduction**: Present the purpose and context of the analysis.
    - **Problem Statement**: Clearly define the problem or question you aim to address.
    - **Results and Discussion**: Showcase your visualizations, insights, and explanations.
    - **Conclusion**: Summarize key findings and any implications.

5. Tools:
- You have the flexibility to choose relevant packages and tools for analysis, with a suggestion to consider **Seaborn** for visualization.

6. Deadline:
- The submission deadline for the Notebook is set for **2024-01-16**.

7. Submission:
- Share your completed **Google Colab Notebook** or **Jupyter Notebook** with me at the email address bernardlkb@ukm.edu.my.

Approach this assignment with the mindset of producing a comprehensive and well-organized document. Your goal is to make the data and its story accessible and understandable to stakeholders.

| Criteria | Marks | | |
| --- | --- | --- | --- |
| Reproducibility | **3**<br>The notebook is<br>100% reproducible | **2**<br>The notebook is<br>reproducible with a few missing<br>steps | **1**<br>The notebook is<br>not reproducible |
| Plots | **10**<br>All the plots are<br>i. suitable,<br>ii. easy to understand<br>iii. observations are properly explained | **7**<br>Some of the plots are<br>i. suitable,<br>ii. easy to understand<br>iii. observations are properly<br>explained | **5**<br>The plots are<br>i. not suitable,<br>ii. hard to understand<br>iii. observations are<br>poorly explained |
| Notebook presentation | **2**<br>The overall notebook is<br>i. properly structured,<br>ii.each section neatly organized,<br>iii. easy to follow | **1**<br>Part of the notebook is<br>i. properly structured,<br>ii.each section neatly organized,<br>iii. easy to follow | **0**<br>The notebook is<br>i. poorly structured,<br>ii. each section is not<br>organized,<br>iii. hard to follow |